

# CSE 351

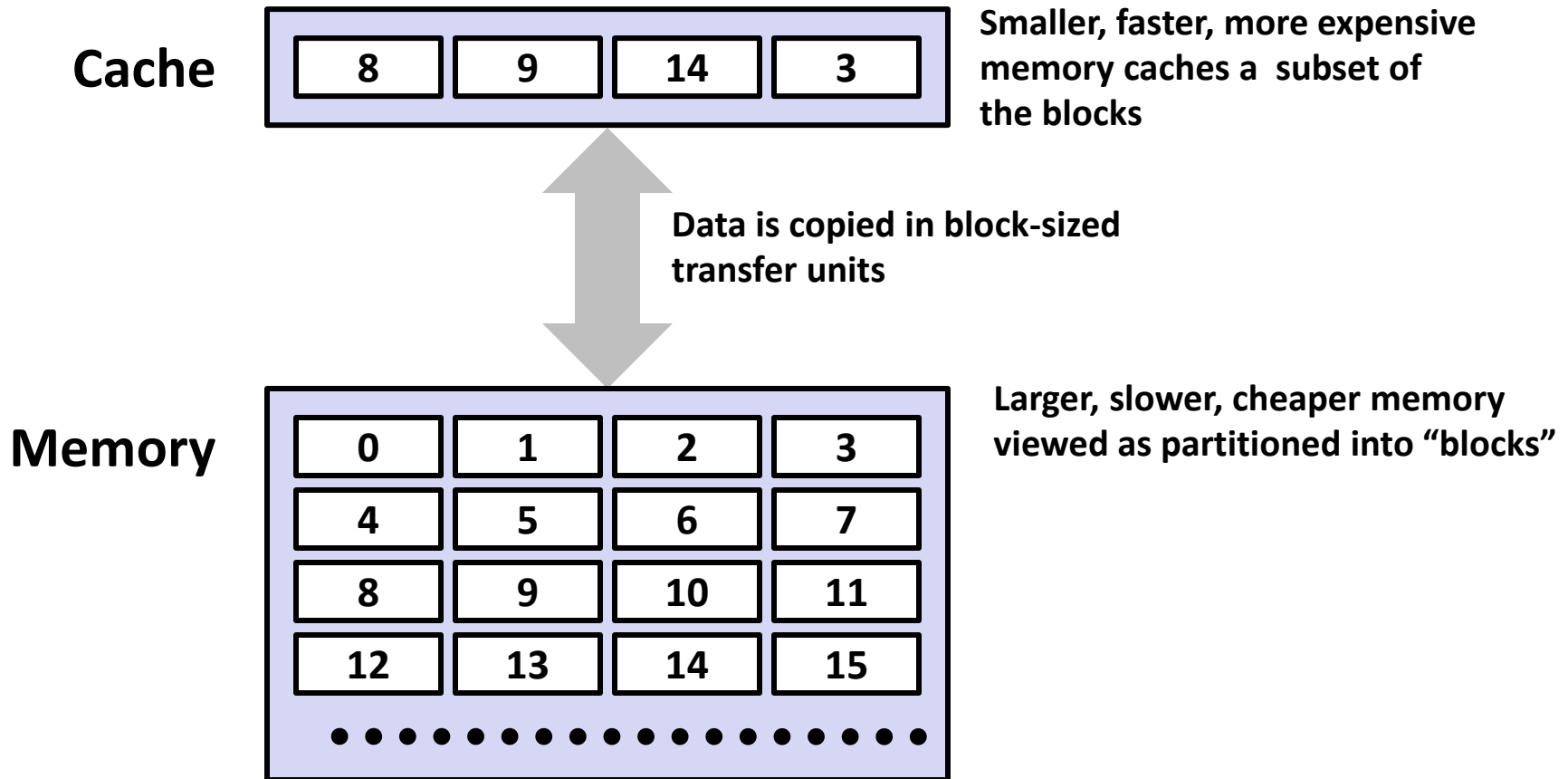
## Section 9

3/1/12

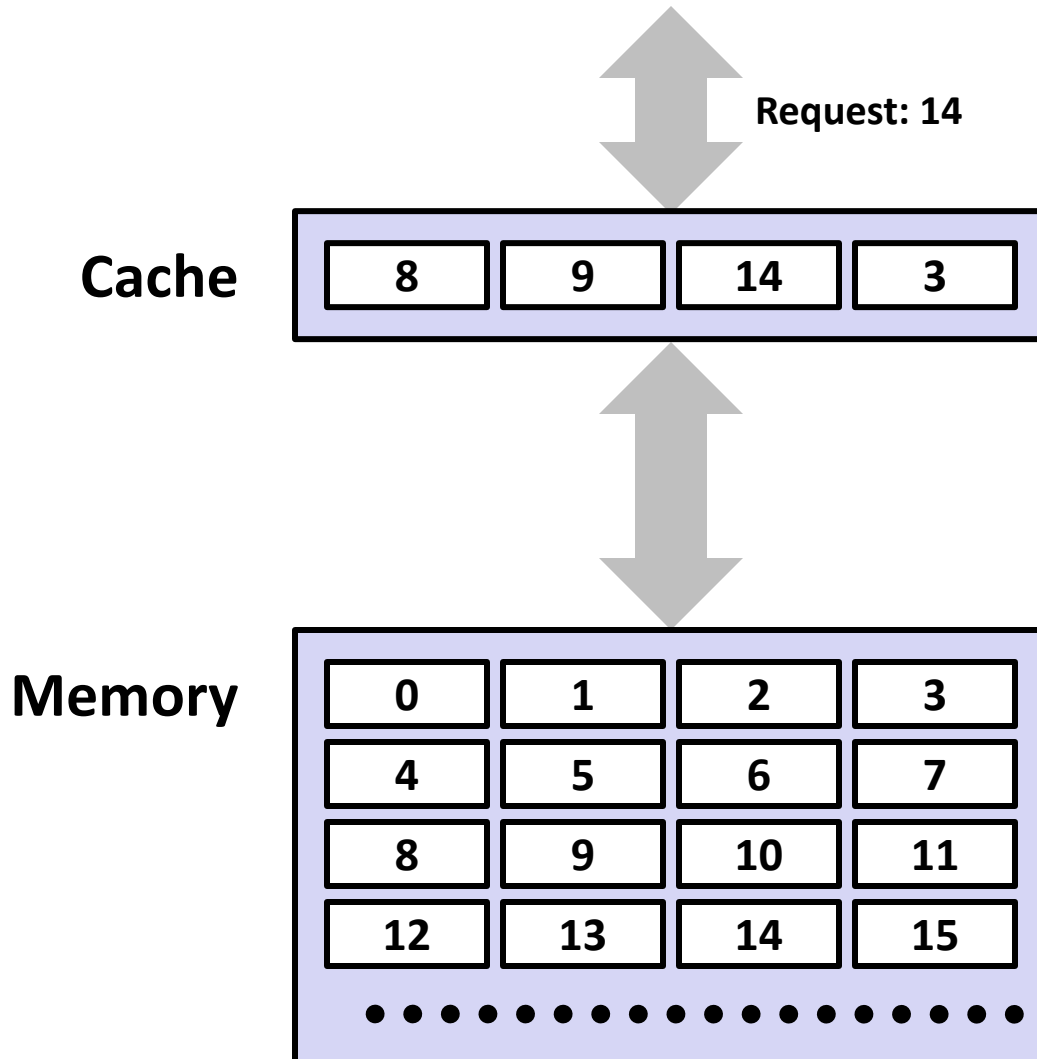
# Agenda

- Caching

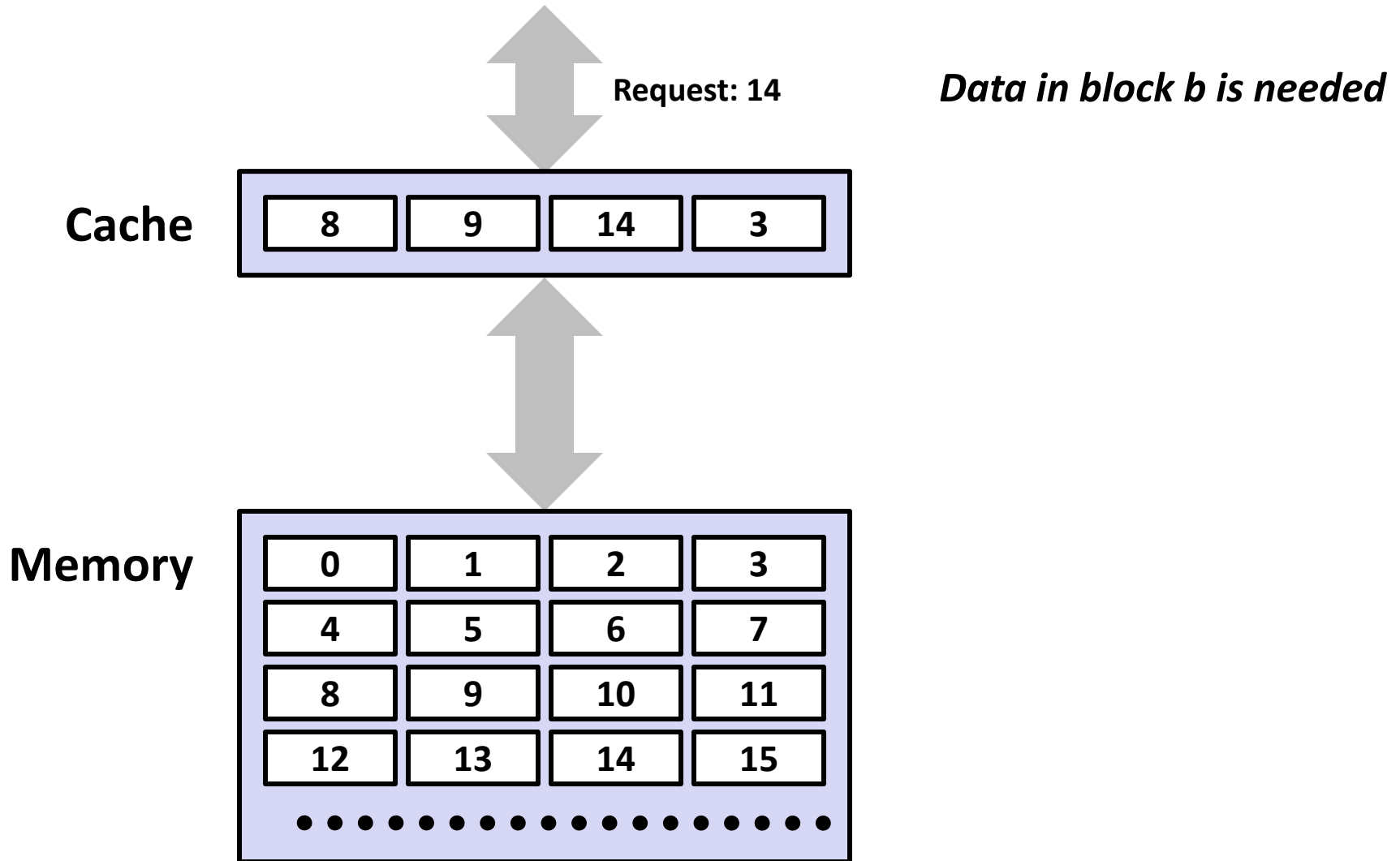
# General Cache Mechanics



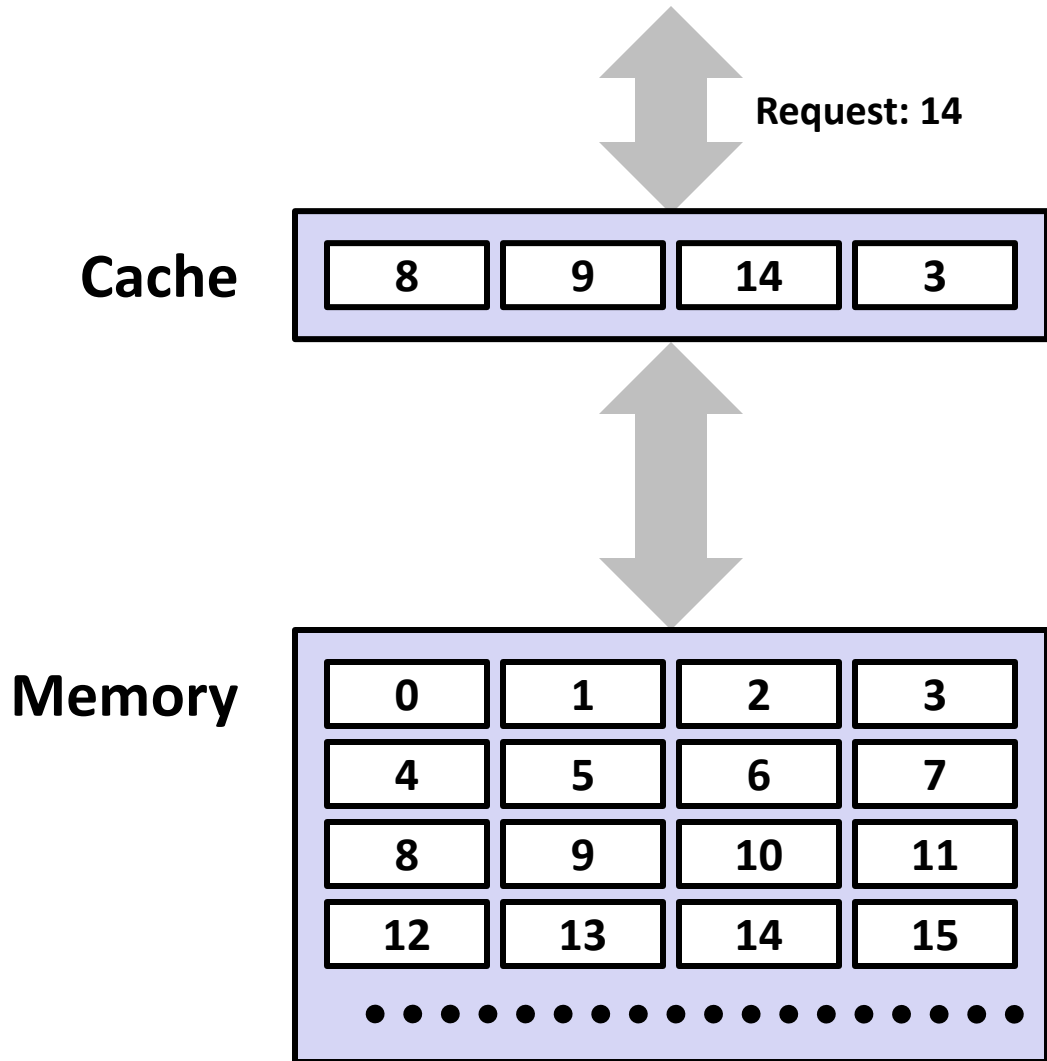
# General Cache Concepts: Hit



# General Cache Concepts: Hit



# General Cache Concepts: Hit

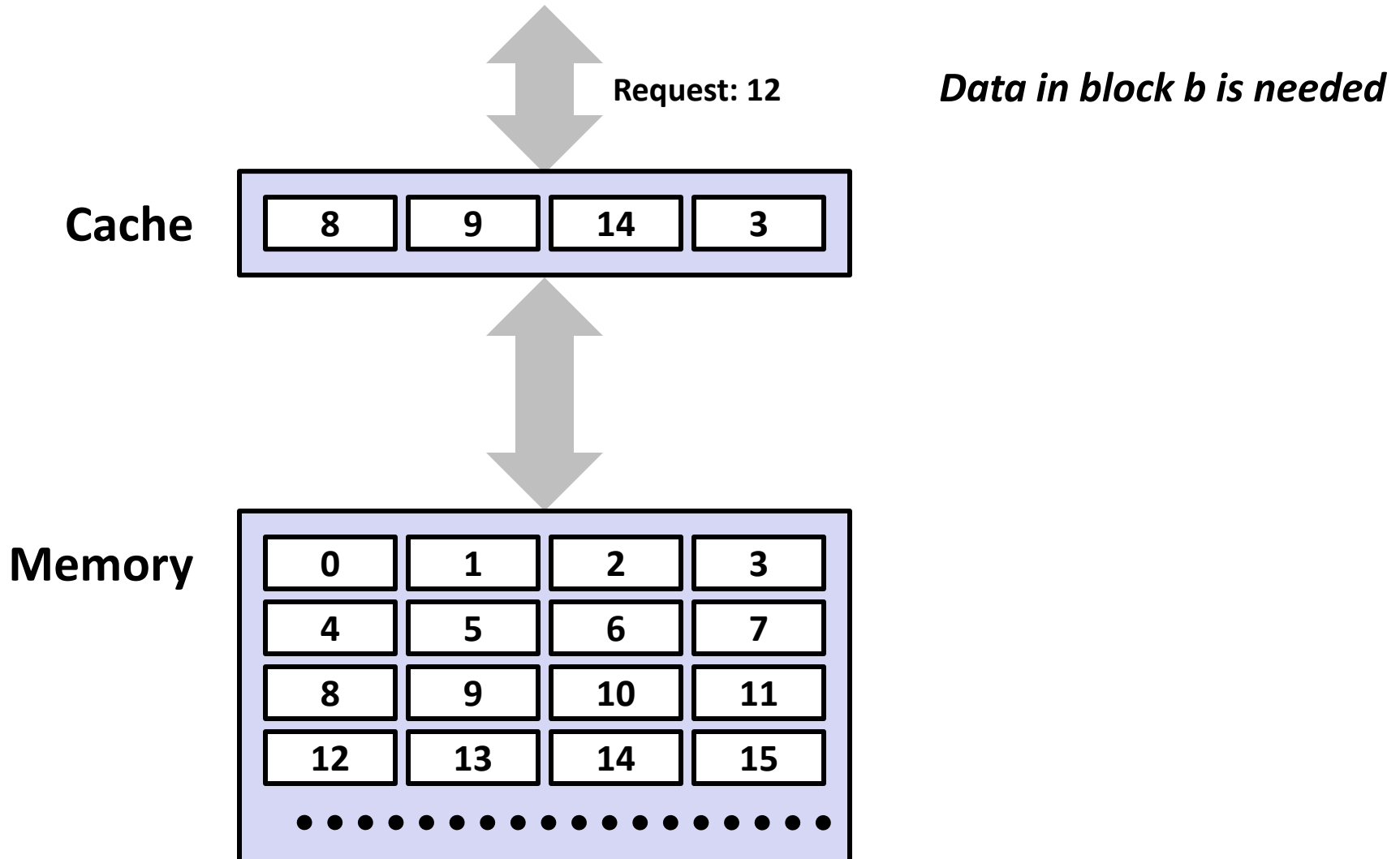


*Data in block b is needed*

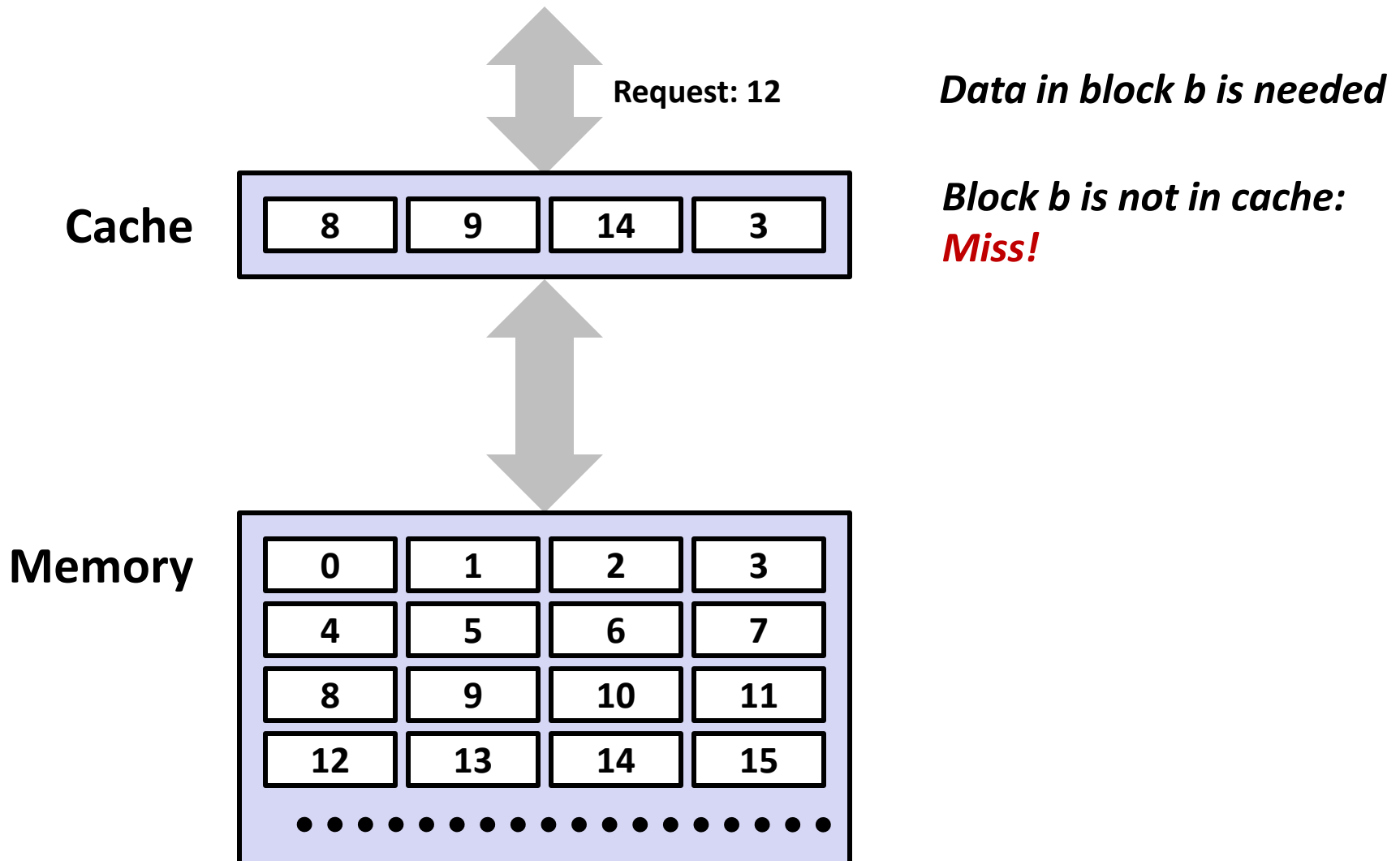
*Block b is in cache:*

***Hit!***

# General Cache Concepts: Miss

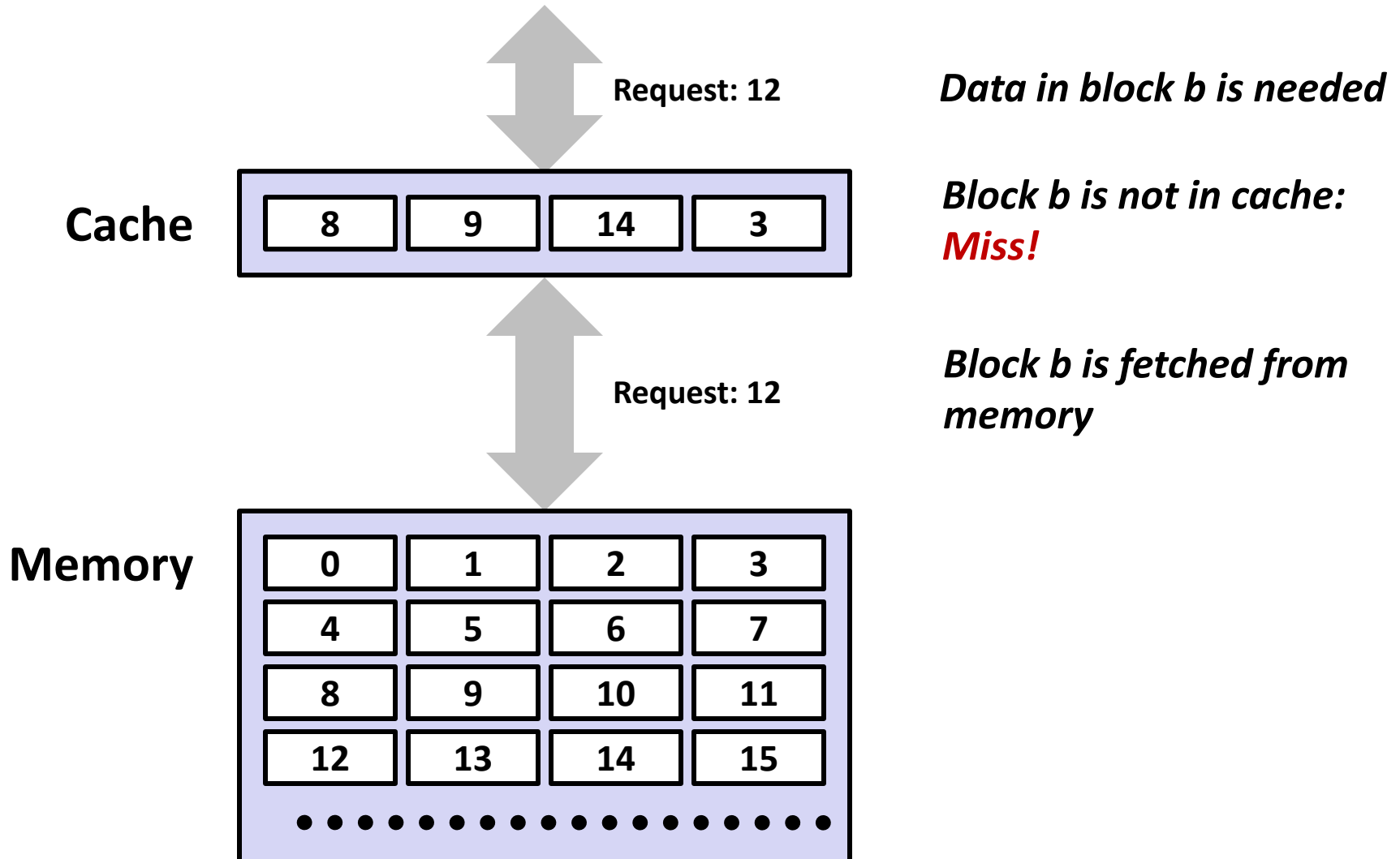


# General Cache Concepts: Miss

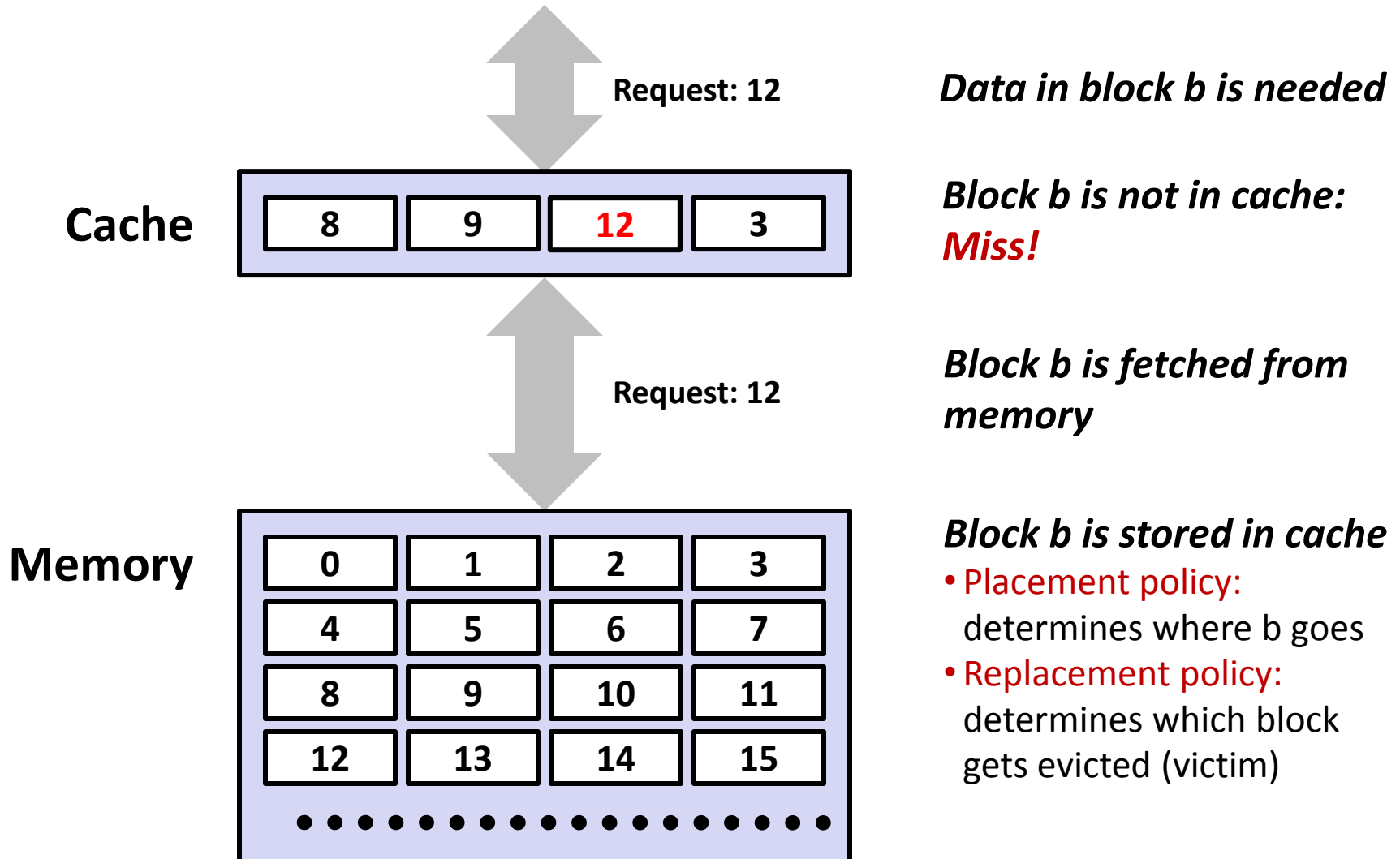




# General Cache Concepts: Miss



# General Cache Concepts: Miss



# Cache Performance Metrics

## ■ Miss Rate

- Fraction of memory references not found in cache (misses / accesses)  
=  $1 - \text{hit rate}$
- Typical numbers (in percentages):
  - 3-10% for L1
  - can be quite small (e.g., < 1%) for L2, depending on size, etc.

## ■ Hit Time

- Time to deliver a line in the cache to the processor
  - includes time to determine whether the line is in the cache
- Typical numbers:
  - 1-2 clock cycle for L1
  - 5-20 clock cycles for L2
  - 30-50 clock cycles for L3

## ■ Miss Penalty

- Additional time required because of a miss
  - typically 100-400 cycles for main memory (**trend: increasing!**)

# Lets think about those numbers

- **Huge difference between a hit and a miss**
  - Could be 100x, if just L1 and main memory
- **Would you believe 99% hits is twice as good as 97%?**
  - Consider:
    - cache hit time of 1 cycle
    - miss penalty of 100 cycles
  - Average access time:
    - 97% hits:  $1 \text{ cycle} + 0.03 * 100 \text{ cycles} = 4 \text{ cycles}$
    - 99% hits:  $1 \text{ cycle} + 0.01 * 100 \text{ cycles} = 2 \text{ cycles}$
- **This is why “miss rate” is used instead of “hit rate”**

# Types of Cache Misses

## ■ Cold (compulsory) miss

- Occurs on first access to a block

## ■ Conflict miss

- Most hardware caches limit blocks to a small subset (sometimes just one) of the available cache slots
  - if one (e.g., block  $i$  must be placed in slot  $(i \bmod \text{size})$ ), direct-mapped
  - if more than one,  $n$ -way set-associative (where  $n$  is a power of 2)
- Conflict misses occur when the cache is large enough, but multiple data objects all map to the same slot
  - e.g., referencing blocks 0, 8, 0, 8, ... would miss every time

## ■ Capacity miss

- Occurs when the set of active cache blocks (the working set) is larger than the cache (just won't fit)

# Why Caches Work

- **Locality:** Programs tend to use data and instructions with addresses near or equal to those they have used recently

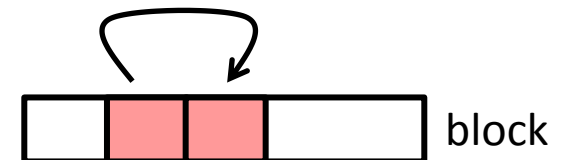
- **Temporal locality:**

- Recently referenced items are *likely* to be referenced again in the near future



- **Spatial locality:**

- Items with nearby addresses *tend* to be referenced close together in time
- How do caches take advantage of this?



# Example: Locality?

```
sum = 0;
for (i = 0; i < n; i++)
    sum += a[i];
return sum;
```

## ■ Data:

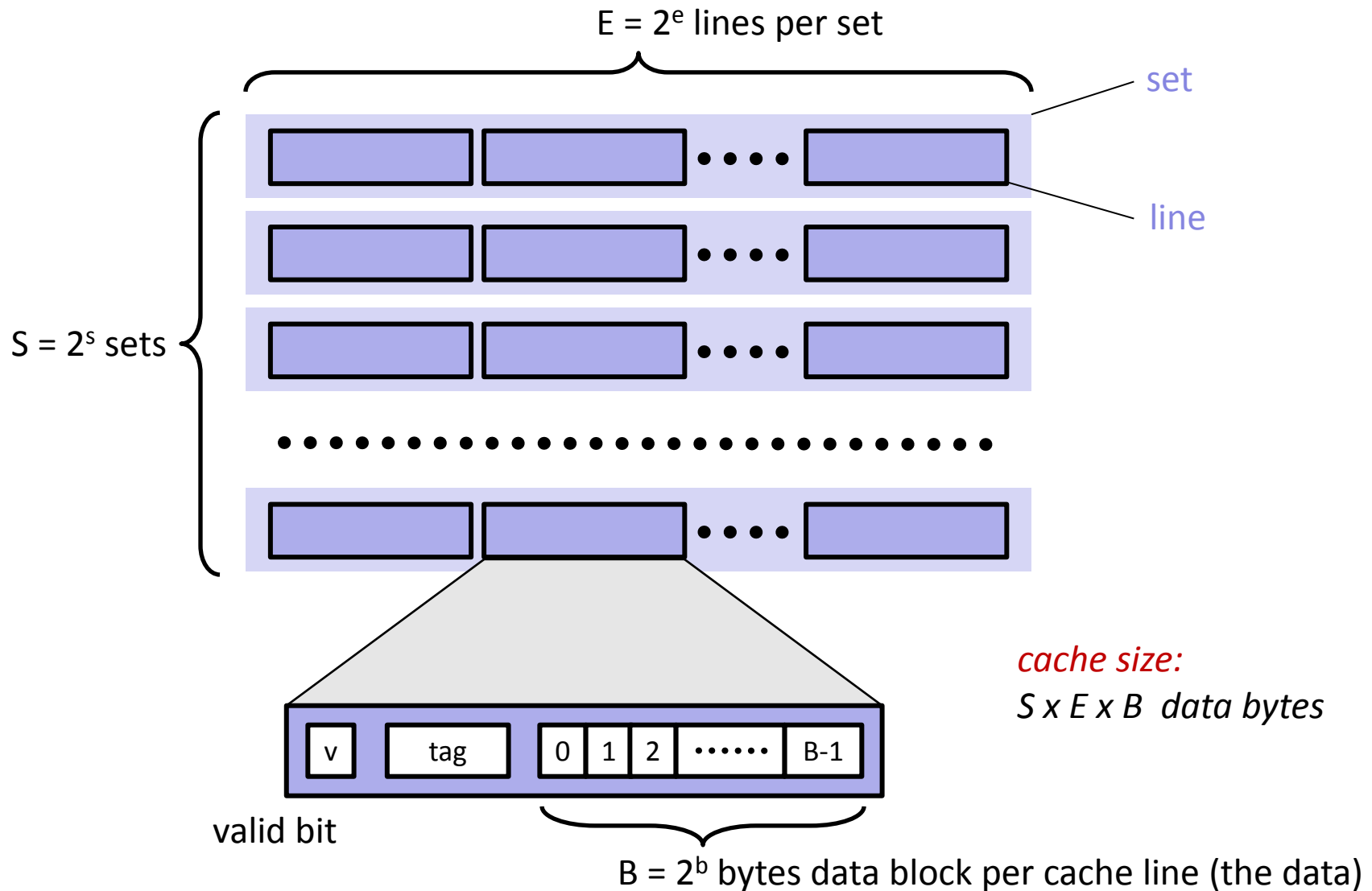
- Temporal: **sum** referenced in each iteration
- Spatial: array **a [ ]** accessed in stride-1 pattern

## ■ Instructions:

- Temporal: cycle through loop repeatedly
- Spatial: reference instructions in sequence

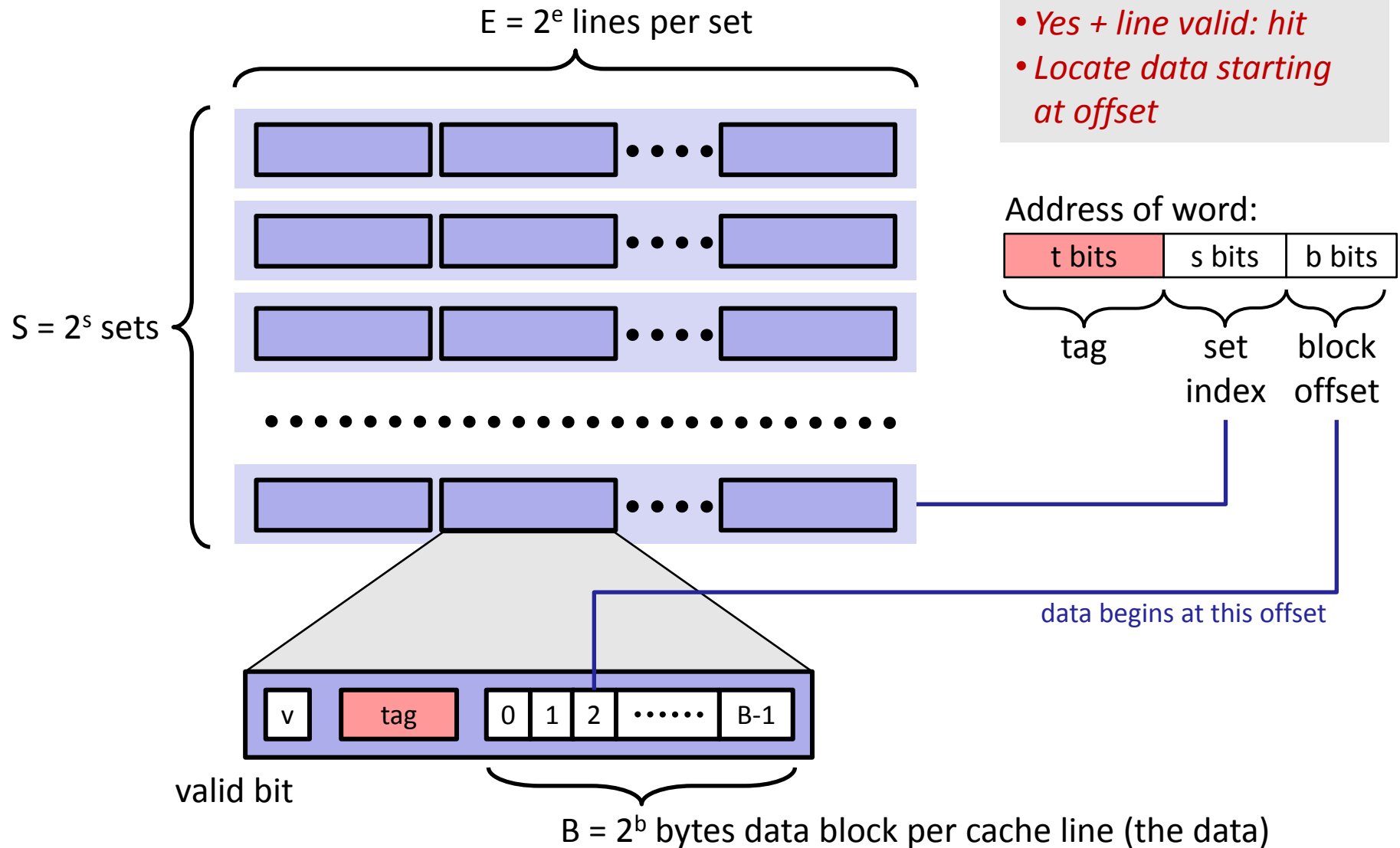
- **Being able to assess the locality of code is a crucial skill for a programmer**

# General Cache Organization (S, E, B)





# Cache Read

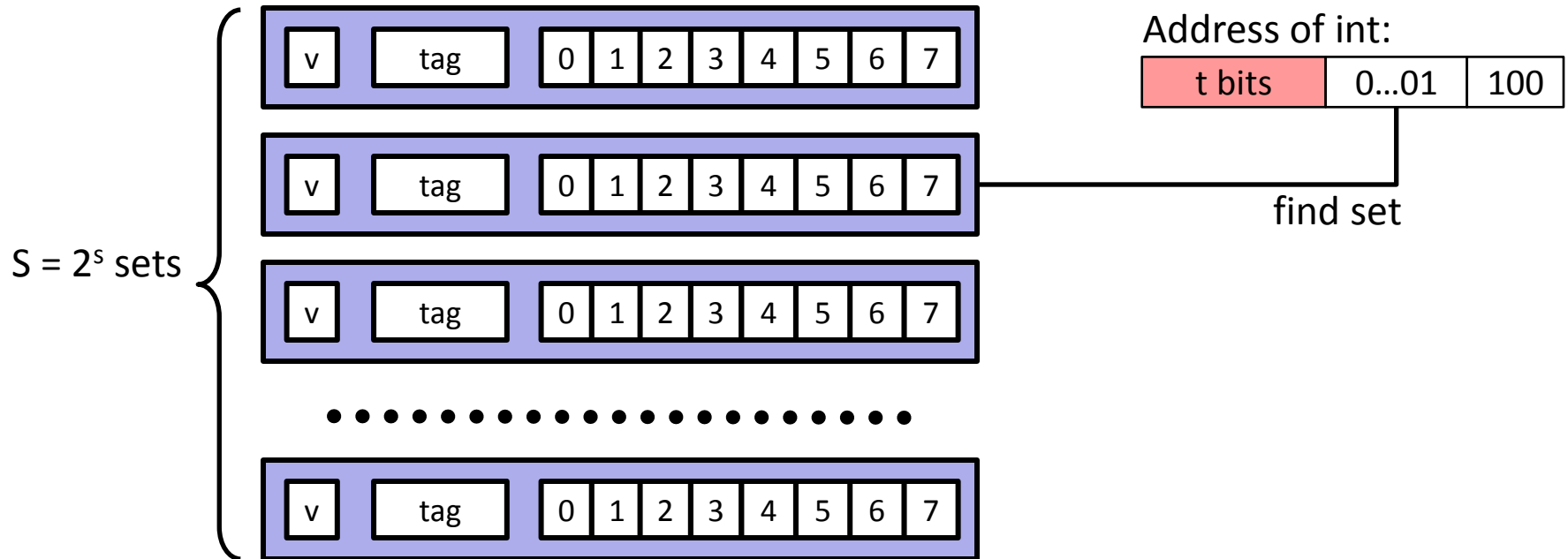


- *Locate set*
- *Check if any line in set has matching tag*
- *Yes + line valid: hit*
- *Locate data starting at offset*

# Example: Direct-Mapped Cache (E = 1)

Direct-mapped: One line per set

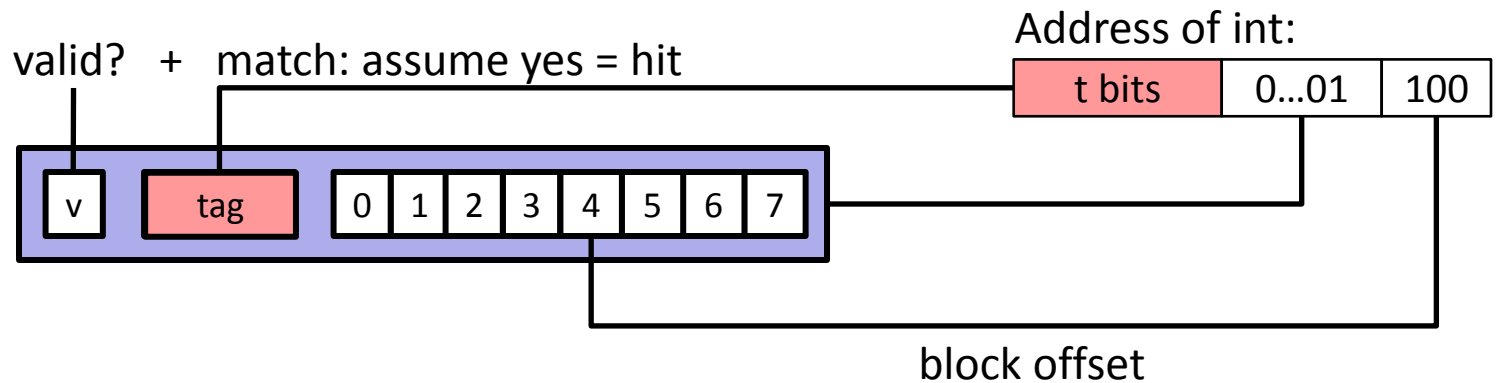
Assume: cache block size 8 bytes



# Example: Direct-Mapped Cache (E = 1)

Direct-mapped: One line per set

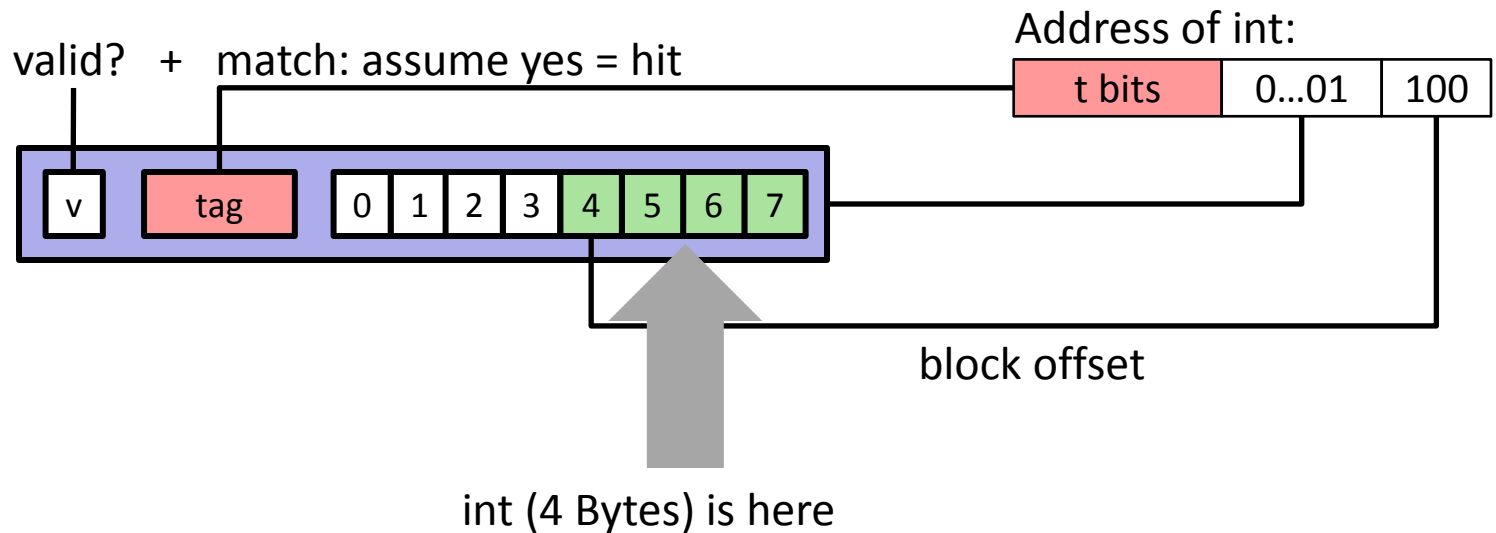
Assume: cache block size 8 bytes



# Example: Direct-Mapped Cache (E = 1)

Direct-mapped: One line per set

Assume: cache block size 8 bytes



**No match:** old line is evicted and replaced

# Example (for E =1)

```
int sum_array_rows(double a[16][16])
{
    int i, j;
    double sum = 0;

    for (i = 0; i < 16; i++)
        for (j = 0; j < 16; j++)
            sum += a[i][j];
    return sum;
}
```

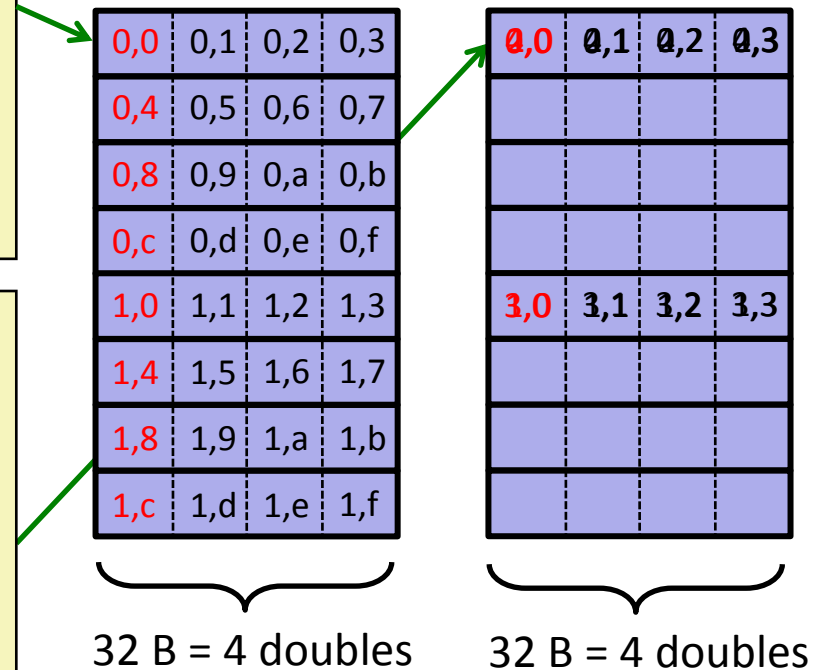
```
int sum_array_cols(double a[16][16])
{
    int i, j;
    double sum = 0;

    for (j = 0; j < 16; j++)
        for (i = 0; i < 16; i++)
            sum += a[i][j];
    return sum;
}
```

Assume *sum, i, j* in registers  
 Address of an aligned element  
 of *a*: *aa...aaaxxxxxyyyy000*

Assume: cold (empty) cache  
 3 bits for set, 5 bits for byte

*aa...aaaxxx xyy yy000*

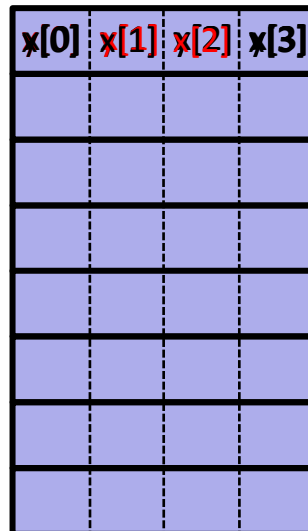


# Example (for E = 1)

```
float dotprod(float x[8], float y[8])
{
    float sum = 0;
    int i;

    for (i = 0; i < 8; i++)
    sum += x[i]*y[i];
    return sum;
}
```

if x and y have aligned  
starting addresses,  
e.g.,  $\&x[0] = 0$ ,  $\&y[0] = 128$



if x and y have unaligned  
starting addresses,  
e.g.,  $\&x[0] = 0$ ,  $\&y[0] = 144$

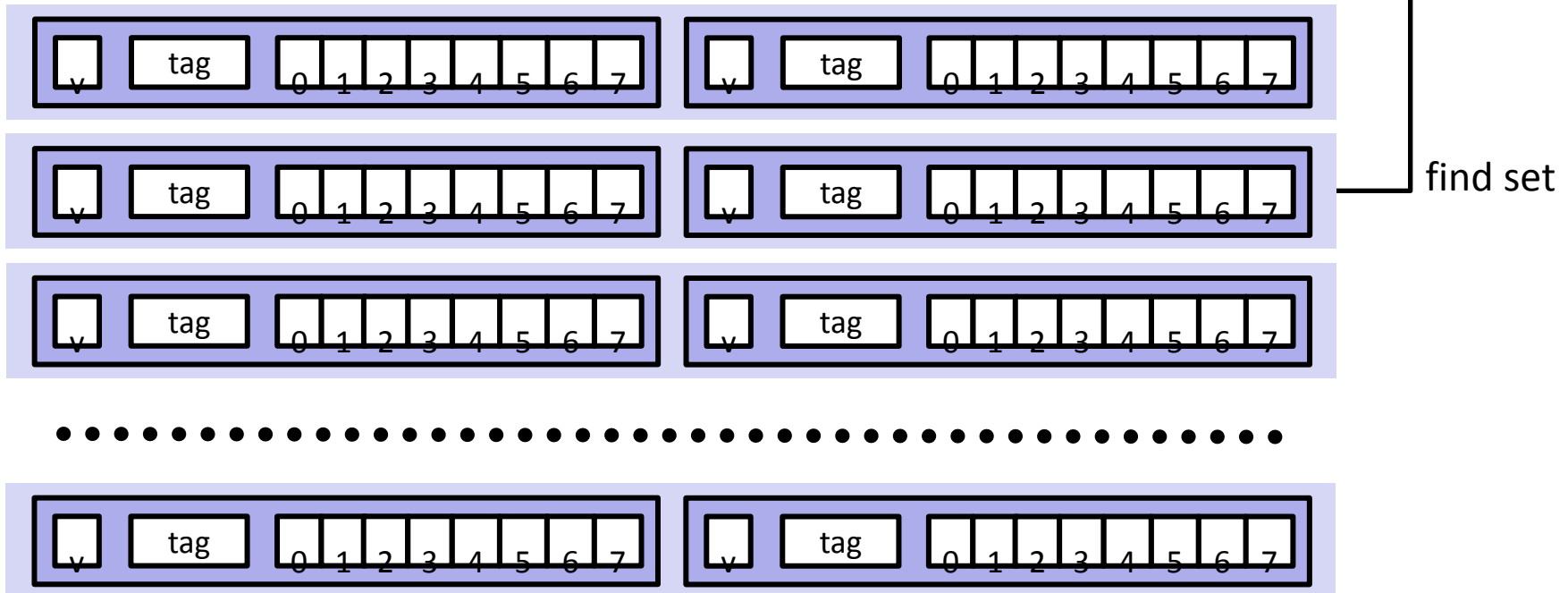
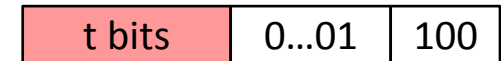


# E-way Set-Associative Cache (Here: E = 2)

E = 2: Two lines per set

Assume: cache block size 8 bytes

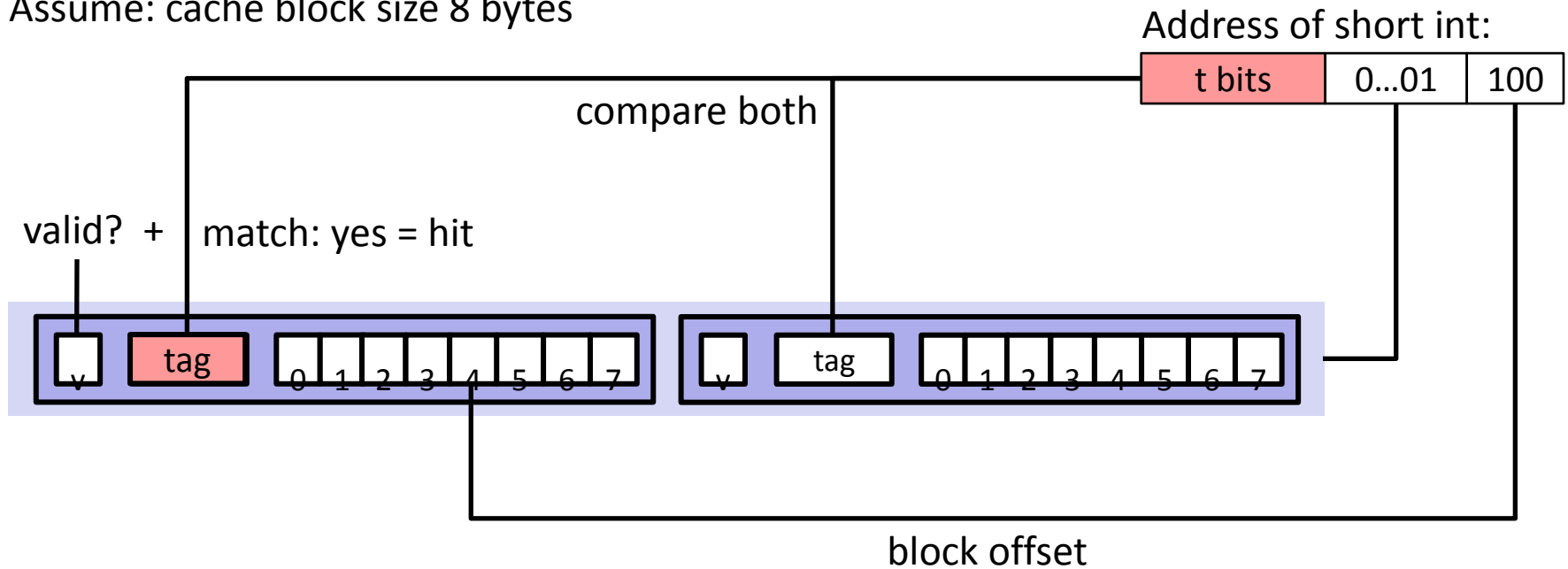
Address of short int:



# E-way Set-Associative Cache (Here: E = 2)

E = 2: Two lines per set

Assume: cache block size 8 bytes

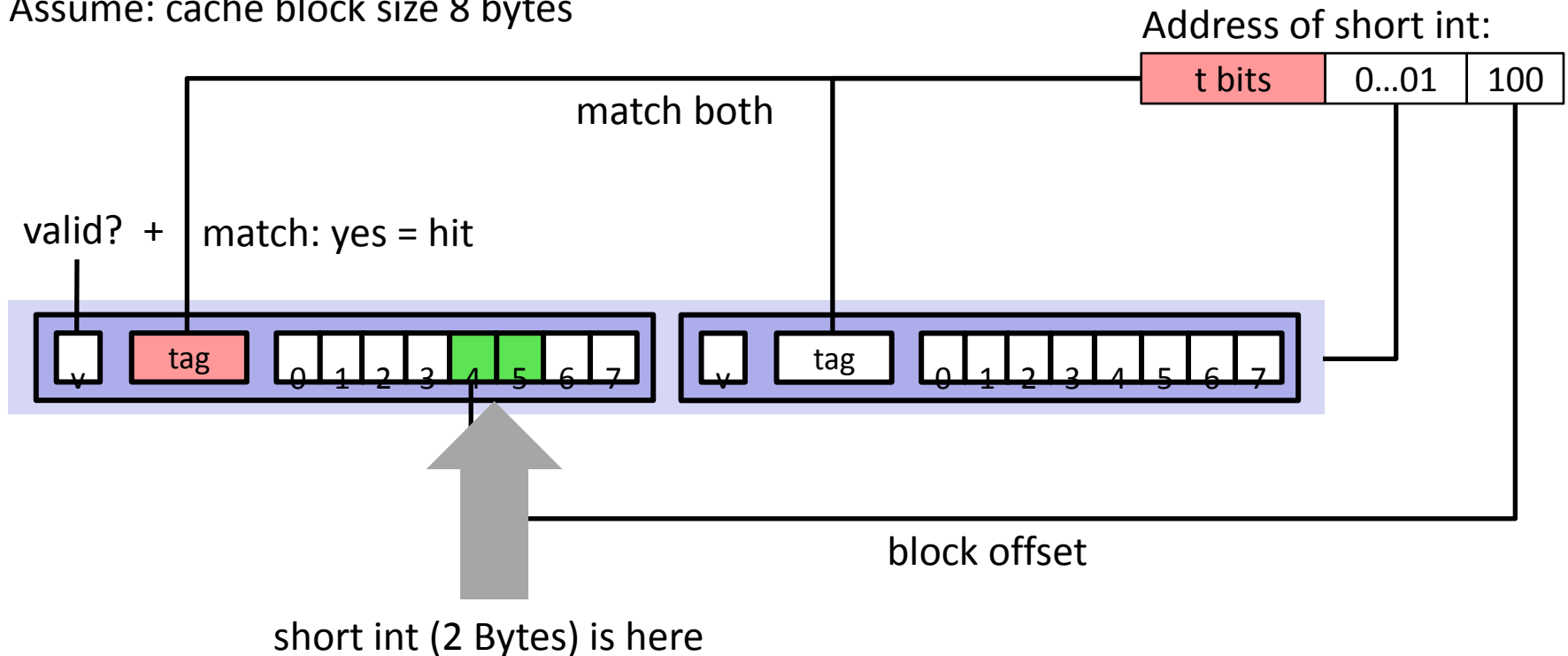




# E-way Set-Associative Cache (Here: E = 2)

E = 2: Two lines per set

Assume: cache block size 8 bytes



## No match:

- One line in set is selected for eviction and replacement
- Replacement policies: random, least recently used (LRU), ...

# Example (for E = 2)

```
float dotprod(float x[8], float y[8])
{
    float sum = 0;
    int i;

    for (i = 0; i < 8; i++)
sum += x[i]*y[i];
    return sum;
}
```

if x and y have aligned  
starting addresses,  
e.g.,  $\&x[0] = 0$ ,  $\&y[0] = 128$   
still can fit both  
because 2 lines in each set

x[0]	x[1]	x[2]	x[3]	y[0]	y[1]	y[2]	y[3]
x[4]	x[5]	x[6]	x[7]	y[4]	y[5]	y[6]	y[7]

# Fully Set-Associative Caches ( $S = 1$ )

- **All lines in one single set,  $S = 1$** 
  - $E = C / B$ , where  $C$  is total cache size
  - $S = 1 = (C / B) / E$
- **Direct-mapped caches have  $E = 1$** 
  - $S = (C / B) / E = C / B$
- **Tags are more expensive in associative caches**
  - Fully-associative cache,  $C / B$  tag comparators
  - Direct-mapped cache, 1 tag comparator
  - In general,  $E$ -way set-associative caches,  $E$  tag comparators
- **Tag size, assuming  $m$  address bits ( $m = 32$  for IA32)**
  - $m - \log_2 S - \log_2 B$

# What about writes?

- **Multiple copies of data exist:**
  - L1, L2, L3, Main Memory, Disk
- **What to do on a write-hit?**
  - Write-through (write immediately to memory)
  - Write-back (defer write to memory until replacement of line)
    - Need a dirty bit (line different from memory or not)
- **What to do on a write-miss?**
  - Write-allocate (load into cache, update line in cache)
    - Good if more writes to the location follow
  - No-write-allocate (writes immediately to memory)
- **Typical**
  - Write-through + No-write-allocate
  - Write-back + Write-allocate